

1 Transcription And Reporting System

2 Barry M. Arons

3 Jeremy Belldina

4 Matthew T. Marx

5 Atty Mullins

6 Haleh Partovi

7 Orion Richardson

8
9 BACKGROUND OF THE INVENTION10 Field of the Invention

11 The present invention relates to transcription and
12 reporting, and specifically to a web-based transcription
13 and reporting tool for use with voice applications.

14
15 Discussion of the Related Art

16 Telephones are ubiquitous in marketplaces around the
17 world. Therefore, many attempts have been made to use the
18 telephone to facilitate electronic commerce. Recent
19 developments in telephone electronic commerce include the
20 use of voice information to guide a transaction between a
21 customer and a voice system. Voice information includes
22 commands spoken by a speaker (e.g. a telephone user),
23 wherein the commands represent transactions between the
24 speaker and the system. For example, commands spoken may
25 include keywords that navigate a menu tree. The spoken
26 commands, called utterances, are interpreted for the voice
27 system by a speech recognizer. Correct interpretation of
28 these utterances by the speech recognizer is key to the
29 success of this method of electronic commerce.

30 In improving the automated interpretation of
31 utterances, voice systems usually use some form of
32 utterance transcription to improve the accuracy of the

1 speech recognizer. Utterances (i.e. audio information) are
2 converted to text information in a process known as
3 transcription. Transcription of utterances allows analysis
4 of the accuracy of the speech recognizer by comparing the
5 result of the speech recognizer to the text information
6 generated by the transcription process. Utterances are
7 typically transcribed with labels, which provide additional
8 information on the utterances. For example, an utterance
9 may be labeled with the gender of a speaker. Different
10 uses for utterances require different labeling schemes.
11 Thus, labels are non-standard over different applications.
12 For example, utterances recorded from a cellular telephone
13 may require labels describing call signal quality.

14 Most transcription and labeling tasks are accomplished
15 with specialized and/or proprietary tools. Such tools
16 range from foot pedal controlled tape players used in
17 conjunction with a typewriter, wherein a transcriptionist
18 listens to the tape and types the results, to custom
19 software that aids in capturing a particular linguistic
20 labeling scheme. Many transcription processes are
21 inefficient in aiding the transcriptionist for both
22 labeling and transcription. For example, in a foot pedal
23 controlled tape player process, a transcriptionist must
24 manually type every utterance and label, thereby having a
25 maximum transcription rate corresponding to the typing
26 speed of the transcriptionist. Additionally, the labels
27 and annotations required for the labeling scheme of the
28 particular application must be remembered or available for
29 reference.

30 Typically, custom software is developed for use with a
31 particular operating system, such as the Macintosh OS,
32 Unix, or Windows NT. The general applicability of such

1 tools is limited by their narrow focus on a specific
2 application, a specific proprietary architecture, or a
3 particular operating system. Due to typically narrow
4 design requirements, custom software is often difficult to
5 extend to differing transcription applications. Moreover,
6 changes to the content and appearance of reports, once
7 initially defined by the custom software, may be limited.
8 Additionally, the requirement of a particular operating
9 system for the custom software limits the flexibility of
10 the transcriptionist in using a particular operating system
11 or associated hardware. Furthermore, some custom software
12 may require on-site transcription, thereby limiting the
13 workforce available for transcription.

14 There are many similar tools for transcription,
15 labeling, and annotation in existence today. Choosing the
16 right combination of tools for a particular application can
17 be a complex decision restricting the later flexibility of
18 the application.

19 Therefore, a need arises for a method of, and a system
20 for, an efficient transcription process having flexible use
21 requirements.

22 23 SUMMARY OF THE INVENTION

24 In accordance with the present invention, a cross-
25 platform transcription and reporting system allows quick
26 transcription of large numbers of utterances and provides
27 analysis of the transcription data in logical reports with
28 linked access to underlying data. The system includes
29 time-saving transcription aids such as buttons defining
30 common noise events and anomalies, thereby allowing a
31 single click to replace numerous typed characters. Labels
32 that are typically consistent across related utterances are

1 pre-defined for each successive related utterance (i.e.
2 consistent labels are "sticky"), thereby obviating the need
3 for the transcriptionist to re-label the related
4 utterances. These transcription aids additionally may be
5 accessed via keyboard shortcuts, thereby saving additional
6 time by allowing a single or multi-key keystroke to replace
7 maneuvering a pointer to click a button and preventing the
8 removal of the transcriptionist's hands from the keys on
9 the keyboard. The text entry box can be pre-loaded with
10 the result of the speech recognizer. In this manner, if
11 the result is correct, the transcriptionist can accept that
12 result by merely hitting the enter key. Note that the text
13 entry box permits only allowable characters, thereby
14 reducing the chance of an incorrect transcription.

15 Features common to web tools such as browsers are
16 taken advantage of in the transcription process, such as
17 auto-completion of a portion of a typed word.
18 Additionally, the use of a web-based system allows
19 distributed transcription across multiple sites and
20 multiple transcriptionists, thereby decreasing costs
21 associated with transcription. For example, multiple
22 transcriptionists, each working from a home location remote
23 from a central database pre-transcribed information, may
24 access the central database simultaneously.

25 Transcribed data are stored in tuples (data
26 structures) along with relevant environment and parameter
27 data. Environment data stored in the tuple includes the
28 grammar-in-use for the utterance. Accordingly, the
29 transcribed data may be compared to the grammar-in-use for
30 in-grammar/out-of-grammar determinations. Additionally,
31 either the audio file of the associated utterance or a
32 pointer to the audio file of the associated utterance is

1 stored in each tuple. Thus, each transcribed utterance may
2 be associated with the original audio utterance.

3 Reports are generated from the tuples meeting a set of
4 reporting criteria. Reports detail the analysis of a set
5 of parameters of the speech recognizer. Reports are
6 presented in one of a set of standard forms, wherein all
7 standard forms include drill-down linking to increasingly
8 detailed levels of supporting data. Because tuples include
9 both the transcribed data and the grammar-in-use, analysis
10 may be made on utterances both in-grammar and out-of-
11 grammar. Accuracy analysis easily includes both mis-
12 accepted results of the speech recognizer and mis-rejected
13 results of the speech recognizer. This ease of generating
14 detailed reports allows authors of a grammar to quickly
15 determine potential grammar issues, such as too large a
16 grammar, too narrow a range of grammar pronunciations, and
17 insufficient limitation of possible utterances.

18 Links to supporting data within the reports allow a
19 double check of the transcription process. For example, a
20 given accuracy statistic, which provides links leading to
21 the audio utterance, allows the audio utterance to be
22 compared to the transcribed utterance. Consistently
23 incorrect results of the speech recognizer indicate an area
24 of training required for the speech recognizer.

25 26 BRIEF DESCRIPTION OF THE DRAWINGS

27 Figure 1 is a block diagram of an utterance storage
28 system in accordance with one embodiment of the present
29 invention.

30 Figure 2 is a screen shot of a sign-in screen for a
31 transcription system according to one embodiment of the
32 present invention.

Figure 3 is a screen shot of a per-call-labels screen according to one embodiment of the present invention.

Figure 4A is a screen shot of a transcription screen according to one embodiment of the present invention.

Figure 4B is another screen shot of the transcription screen of Figure 4A according to one embodiment of the present invention.

Figure 5 is a flow diagram of the transcription process according to one embodiment of the present invention.

Figure 6 is a screen shot of a top-level drill-down report according to one embodiment of the present invention.

Figure 7 is a screen shot of a first-level-down drill-down report according to one embodiment of the present invention.

Figure 8 is a screen shot of a second-level-down drill-down report according to one embodiment of the present invention.

Figure 9 is a screen shot of a third-level-down drill-down report according to one embodiment of the present invention.

Similar elements in the above Figures are labeled similarly.

DETAILED DESCRIPTION OF THE DRAWINGS

In accordance with the present invention, a cross-platform transcription and reporting system provides ease of use and user access from multiple locations. Web-based transcription tools allow multiple transcriptionists to interface with the information database using a web browser. Transcription information is compiled in a

1 variety of reports organized in a drill-down to detail
2 fashion. Specifically, direct access is provided from top-
3 level statistics to low-level detail through a series of
4 hyperlinks. A hyperlink (link) is an element in a web page
5 that, when clicked upon, provides access to another web
6 page, typically by navigating the web browser to the other
7 web page. Web-based transcription tools additionally allow
8 the use of built-in browser features (e.g. the auto-
9 complete function).

10 In a telephone-based speech recognition system, during
11 a transaction, users are led through a series of voice
12 menus to achieve a desired result. For example, a
13 transaction may include the user choosing a first voice
14 option from a main menu (e.g. information regarding
15 "weather"), and a second voice option from a secondary menu
16 (e.g. desired location of weather information is "San Jose,
17 California"). To increase the accuracy of the speech
18 recognition system, each menu has an associated local
19 grammar with a limited scope. A grammar defines the set of
20 valid expressions that a user can say when interacting with
21 the speech recognition system. For example, a local
22 grammar for the main menu above may include the expressions
23 "stock quotes", "traffic", and "weather". A local grammar
24 for the weather secondary menu may include the expressions
25 "Chicago, Illinois", "New York City, New York", and "San
26 Jose, California". To limit the scope of the local grammar
27 in the secondary menu, the expressions "stock quotes",
28 "traffic", and "weather" from the local grammar in the
29 primary menu are not valid expressions when interacting
30 with the secondary menu. Thus, the main menu local grammar
31 is not in use when interacting with the secondary menu.
32 Note that menus may have multiple associated local

1 grammars. For example, the secondary menu above may also
2 have additional local grammars, such as a list of valid zip
3 codes corresponding to the city/state pairs of the first
4 local grammar.

5 Intrinsic grammars are also available for use with
6 menus. Intrinsic grammars are grammars with widespread
7 applicability. Some intrinsic grammars are always
8 available and may be used at any time when interacting with
9 menus. For example, a global commands intrinsic grammar
10 may include the expressions "help", "go back", and
11 "repeat". In one embodiment, because these global commands
12 are useful for all menus, the global commands intrinsic
13 grammar is always available. Other intrinsic grammars,
14 such as a telephone number grammar (recognizing strings of
15 numbers), and a date/time grammar (recognizing days of the
16 week, months, days, and years) are available for use with
17 appropriate menus.

18 Utterances from a telephone-based speech recognition
19 system are recorded and used to train the speech
20 recognition system. Utterances are the sounds made by a
21 user (speaker) of the speech recognition system.
22 Recordings of these utterances (e.g. typically 1 to 5
23 seconds) are digitized and stored in a database or a file
24 system hierarchy (database). This database consists of
25 both the utterance recordings (utterances) and a log of
26 information relating to those utterances (such as the time
27 the utterance was made, the grammar then in use, the result
28 of the speech recognizer, other parameters, and a pointer
29 to the specific utterance recording). Each stored element
30 may be described as a record tuple: a series of records,
31 each record having multiple elements. In one embodiment,
32 each record is listed in the form (date/time, grammar then

1 in use, result, parameters, pointer to stored utterance
2 recording). In one embodiment, the utterance recording
3 replaces the pointer to stored utterance recording in the
4 tuple.

5 Figure 1 is a block diagram of an utterance storage
6 system 100 in accordance with one embodiment of the present
7 invention. Storage system 100 includes hosting sites 101
8 and 102, which are physical locations housing storage
9 equipment. Each hosting site includes one or more pods
10 (e.g. hosting site 101 includes pods 105 and 106, and
11 hosting site 102 includes pod 107). A pod is a collection
12 of telephony speech recognition equipment coupled to phone
13 lines. Each pod can handle a given number of simultaneous
14 users (callers) interfacing with the speech recognition
15 system. Thus, each pod creates utterance recordings from
16 the user and generates a log file containing the associated
17 record tuples.

18 Due to the volume of data (i.e. utterance recordings
19 and the log file) stored in pods 105-107, selection
20 criteria can be applied by filters 108 and 109 to aggregate
21 the data from pods 105-107 into one or more tiers of
22 intermediate storage 103 and 104. For example, in one
23 embodiment, filter 108 applies selection criteria to the
24 data in pods 105-107 to retrieve 50% of the data in pods
25 105-107 each evening and store that data in intermediate
26 storage 103. In this embodiment, filter 109 applies
27 selection criteria to the data retrieved through the use of
28 filter 108, such as removing data attributable to internal
29 callers (internal users) testing the speech recognition
30 system. In this way, data to be transcribed can be
31 filtered prior to transcription into meaningful groups with
32 associated general characteristics for later transcription.

1 Once the data has been created and filtered, the
2 transcription process begins. Because the present cross-
3 platform transcription system is web-based,
4 transcriptionists may transcribe data from any location
5 having a suitable connection to the data. Data may be
6 accessed over a network using an Internet protocol, such as
7 hypertext Transfer Protocol (HTTP). HTTP is an
8 application-level protocol for distributed, collaborative,
9 hypermedia information systems. In one embodiment, the
10 network used is a Virtual Private Network (VPN). A VPN
11 uses privacy features such as encryption and tunneling to
12 connect users or sites over a public network, typically the
13 Internet. In comparison, a private network uses dedicated
14 lines between each point on the network. As described in
15 more detail below, a transcriptionist first initiates a
16 connection to the database through a web browser, signs
17 into the transcription system, chooses the records to be
18 transcribed, and then begins the transcription process.

19 Figure 2 is a screen shot of a sign-in screen for a
20 transcription system according to one embodiment of the
21 present invention. Web browser 200 (e.g. the Internet
22 Explorer® web browser) displays the address (i.e. location)
23 of the transcription system in address window 201. Web
24 browser 200 displays sign-in screen 200(A). Within sign-in
25 screen 200(A), the transcriptionist chooses a date of files
26 to transcribe (field 205), enters a unique transcriptionist
27 ID (field 206), enters a record starting number (field
28 207), and submits the above information by pressing submit
29 button 210. A record is a collection of utterances during
30 one interface with the speech recognition system (i.e.
31 during one call). Comments may be sent to the system
32 administrators by pressing comment button 211, and a

1 tutorial describing the transcription system may be reached
2 by clicking on tutorial hyperlink 212. In one embodiment,
3 comments may also be stored with the transcribed
4 utterances. Pressing submit button 210 causes the per-
5 call-labels screen (Figure 3) to appear within the web
6 browser window.

7 Some embodiments may offer more sophisticated
8 utterance selection mechanisms in conjunction with sign in
9 to support more selective transcription in response to
10 specific needs. For example, if "driving directions" was
11 introduced as a new application, it might be possible to
12 easily select only "driving direction"-related utterances
13 for transcription. In other embodiments, the
14 transcriptionist may not be directly presented with the
15 utterance selection options, e.g., they may be
16 predetermined for a transcriptions based on her/his login.
17 In this embodiment, one or more supervisors and/or
18 automated processes might automatically select utterances
19 for a particular transcriptionist according to one or more
20 criteria. Also, as will become clearer when discussed
21 below, typically most, or all, of the available utterances
22 for a particular call are transcribed in a single session
23 by a single transcriber. This maximizes the value of the
24 transcriber's natural language capabilities (especially if
25 the transcriber is familiar with the application) and
26 increases accuracy. However, this is not a technical
27 requirement.

28 Figure 3 is a screen shot of a per-call-labels screen
29 according to one embodiment of the present invention. Web
30 browser 200 navigates the browser window to the address
31 shown in address window 201 in response to pressing submit
32 button 210 in sign-in screen 200(A). Thus, a per-call-

1 labels screen 200(B) is shown subsequent to sign-in screen
2 200(A), but prior to each record being transcribed. A
3 series of utterances made during one call to the speech
4 recognition system are likely to share certain
5 characteristics: gender of user, whether user is a native
6 or non-native speaker, car background noise, and overall
7 bad audio quality. By allowing entry of these consistent
8 labels once, labels that are typically consistent
9 throughout a call need only be entered once. As described
10 below, these per-call-labels are then filled in to the
11 transcription screen for each related utterance to be
12 transcribed (i.e. consistent labels are "sticky"), thereby
13 speeding the transcription of each utterance.

14 In one embodiment, the short recording of the first
15 utterance assigned to the first record is automatically
16 played upon initial display of per-call-labels screen
17 200(B). Audio control panel 310 allows the
18 transcriptionist to play the utterance, as well as perform
19 other audio operations such as change the volume and pause
20 the replay of the recording. Once the transcriptionist
21 hears the utterance, per-call-labels 301-304 may be
22 defined. Thus, the user's gender (either male or female)
23 is defined using gender radio button 301 and the user's
24 accent (either native or non-native) is defined using
25 accent radio button 302. A radio button is a device that
26 allows the selection of only one of a group of options
27 (e.g. only one of "male" or "female" may be chosen in radio
28 button 301). Similarly, noise within a car while a user is
29 speaking on a cellular telephone may be noted by checking
30 car noise checkbox 303 and bad audio signal may be noted by
31 checking bad audio checkbox 304. A checkbox is a toggle
32 device that allows a value to be set on (box is checked) or

1 off (box is unchecked). Thus, an unchecked box indicates
2 that the associated attribute is not present in the current
3 utterance (or record). Note that keyboard shortcuts (hot
4 keys) are available for radio buttons 301 and 302 as well
5 as for checkboxes 303 and 304.

6 Note that transcriptionists make educated estimates
7 for some of these values. For example, a transcriptionist
8 may identify a particular utterance with a "female" label
9 by using radio button 301. This transcription label does
10 not mean that the user was in fact a woman, but rather
11 means that the transcriptionist believes the caller to be a
12 female. Throughout the transcription process as described
13 below, the per-call labels may be adjusted as appropriate.

14 Similarly to sign-in screen 200(A), comments may be
15 entered by pressing comment button 211, and a tutorial
16 describing the transcription system may be reached by
17 clicking on tutorial hyperlink 212. Additionally, help on
18 labels may be reached by clicking on "help: labels"
19 hyperlink 312. Pressing submit button 210 causes the
20 transcription system to accept the per-call-labels
21 information and then causes the transcription screen
22 (Figure 4A) to appear within the web browser window.

23 Figure 4A is a screen shot of a transcription screen
24 according to one embodiment of the present invention. Web
25 browser 200 navigates a browser window to the address shown
26 in address window 201 in response to pressing submit button
27 210 in per-call-labels screen 200(B). A transcription
28 screen similar to transcription screen 200(C) is shown for
29 each utterance in a record.

30 The short recording of the utterance to be transcribed
31 is automatically played upon display of transcription
32 screen 200(C). Text entry field 409 is automatically

1 populated with the result of the speech recognizer. If the
2 result of the speech recognizer is correct and no
3 additional labels need be defined, the transcriptionist
4 need only hit "Enter" on the keyboard (the keyboard short
5 cut for submit button 210) to accept the transcription and
6 move onto the next utterance to be transcribed. If the
7 transcriptionist disagrees with the automatically populated
8 text, the transcriptionist types the text translation of
9 the utterance into text entry field 409 in place of the
10 automatically populated text and adds any needed labels.
11 Text entry field 409 is discussed in more detail with
12 respect to Figure 4B below. Note that previous button 421
13 allows the transcriptionist to return to the transcription
14 screens of previously transcribed utterances.

15 In addition to transcribing the utterance, the
16 transcriptionist provides labels describing the utterance
17 sound recording. Per-call-labels 301-304, which were pre-
18 populated from information from per-call-labels screen
19 200(B), are available for alteration in transcription
20 screen 200(C). During one call, a first user may hand the
21 telephone to a second user of a different gender or accent,
22 necessitating a change in one of these "sticky" fields or
23 the first user may move from a house to a car, etc.
24 Additionally, checkboxes are provided for noting such
25 events as background noise during the utterance (background
26 noise checkbox 401) and whether the utterance recording is
27 truncated either at the beginning (beginning cut off
28 checkbox 402A) or at the end (end cut off checkbox 402B).

29 Noise events buttons 410-415 generate labeling text
30 denoting an utterance directed other than towards the
31 speech recognizer (side speech button 410), breath noise
32 (breath noise button 411), a word fragment (fragment button

1 412), a DTMF touchtone noise (touchtone button 413), the
2 sound of a hang up (hang up button 414), or other noise
3 (other noise button 415). For example, pressing side
4 speech button 410 generates the label "[side_speech]" and
5 then inserts that label into text entry box 409 (not
6 shown). Help is available for these noise events by
7 clicking on "help: noise events" hyperlink 405.

8 Anomalies buttons 416-420 insert labeling text into
9 text entry box 409 denoting anomalous utterances, including
10 unintelligible utterances (unintelligible button 416),
11 interjections such as "ah", "uh", or "oh" (ah, uh, oh
12 button 417), and filler noises such as "um", "hmm", and
13 "hum" (um, hmm, hum button 418). Anomalous utterances also
14 include those transcriptions which are the best guess of
15 the transcriptionist (best guess button 419) and which are
16 the correct spelling of a mispronounced word (mispronounced
17 button 420). For example, pressing mispronounced button
18 420 encases the transcribed word in asterisks within text
19 entry box 409 (not shown). Although labels for anomalies
20 are typically nonstandard across transcription systems, the
21 consistent use of one type of label for each type of
22 anomaly allows the possibility of a global label
23 replacement to meet the requirements of a particular
24 reporting system or analysis framework. Help is available
25 for these anomalous utterances by clicking on "help:
26 anomalies" hyperlink 406. Help is available for these
27 transcription conventions by clicking on "help:
28 transcription conventions" hyperlink 407. Note that most
29 buttons, radio buttons, and checkboxes have keyboard
30 shortcuts, thereby allowing the transcriptionist to perform
31 most transcription functions without moving hands away from
32 the keyboard.

Figure 4B is another screen shot of the transcription screen 200(C) according to one embodiment of the present invention. As described above, text entry field 409 is pre-populated with the result of the speech recognizer. If the transcriptionist disagrees with the automatically populated text, the transcriptionist types the text translation of the utterance into text entry field 409 in place of the automatically populated text. As the transcriptionist types in text entry field 409, drop-down selection menu 409A (a part of text entry field 409) appears containing a list of possible words typed by the transcriptionist. As shown, the typed letters "t-e-l-l" produces a list of words beginning with those letters, such as "tell me" and "tell me more". The auto-complete function of web browser 200 may be used to auto-complete the text typed by the transcriptionist with the most frequently used word having the same root letters. Note that drop-down selection menu 409A obscures audio tool 310, play button 311, jump button 421, and a portion of submit button 210 from view within transcription screen 200(C) (see Figure 4A). Once a word is chosen for text entry box 409, drop-down selection menu 409A disappears. In one embodiment, only predetermined characters are allowable. In this embodiment, inserting a character not allowed (e.g. illegal punctuation or a numerical digit) in text box 409 triggers a warning to the transcriptionist that the character is not allowed for the transcription scheme.

Additionally, if supported by web browser 200, the transcriptionist may tab to select each element in turn (e.g. side speech button 410, then breath noise button 411). The transcriptionist may hit the "Enter" key on the keyboard as a short cut to perform the action associated

1 with the highlighted element, thereby allowing the
2 transcriptionist to additionally access most displayed
3 elements without removing hands from the keyboard.

4 Figure 5 is a flow diagram of the transcription
5 process according to one embodiment of the present
6 invention. As described above, a web browser is navigated
7 to the address of the transcription system in step 501.
8 Each transcriptionist signs into the transcription system
9 in step 502 and chooses a starting record number. Steps
10 502 and 503 are performed using sign-in screen 200(A)
11 (Figure 2). Other embodiments of the transcription process
12 include additional steps, such as a transcriptionist
13 verification screen, wherein each transcriptionist verifies
14 authorized access (e.g. uses a password to sign into the
15 transcription system). As noted above with respect to
16 Figure 2, a transcriptionist may be transcribing a subset
17 of a call, e.g., all utterances in "driving directions",
18 etc. However, for convenience the term "call" will be used
19 since in the preferred embodiment, a transcriptionist only
20 works on utterances taken from a single phone call at a
21 time.

22 Additionally, in one embodiment, the utterances from a
23 given call are transcribed in sequence. Because calls
24 navigate through a defined set of menus with defined
25 grammars, transcribing the calls in sequence gives the
26 transcriptionist additional context, thereby improving the
27 transcription accuracy. For example, an utterance such as
28 "San Jose, California" might be difficult to recognize out
29 of context, but may be easier to recognize if the previous
30 utterance was "weather", thereby indicating the desire to
31 obtain weather information including the forecast for a
32 particular city.

1 Once a starting record is chosen in step 503, per-
2 call-labels are defined in step 504 using per-call-labels
3 screen 200(B) (Figure 3). The first utterance is
4 transcribed in step 505 using transcription screen 200(C).
5 If additional utterances are present in the record (step
6 506), the additional utterances are transcribed returning
7 to transcribe utterance step 505. If no more utterances
8 are present in the record, a decision is made by the
9 transcriptionist whether or not to continue transcribing
10 records in step 507. In one embodiment, the
11 transcriptionist initially chooses a certain number of
12 records to transcribe, thereby automating "continue
13 transcription?" step 507.

14 If the transcription is to continue with another
15 record in step 507, the next record is selected in step 508
16 and per-call-labels defined for that record in step 504.
17 If the transcription is finished, the transcription system
18 is exited in step 509.

19 In one embodiment, the transcribed information extends
20 the tuple stored in the database to include an additional
21 data element indicating the transcribed value. For
22 example, after transcription, the tuple contains the
23 elements (date/time, grammar then in use, result,
24 parameters, pointer to stored utterance recording,
25 transcribed result).

26 It is important for all of this transcription data to
27 be available for analysis in a meaningful, yet easy to
28 understand fashion. Accordingly, the present invention
29 provides for a system of drill-down reports to describe the
30 transcription data. These drill-down reports include data
31 compilation into a top-level analysis with direct
32 hyperlinked access to supporting data. As described below,

1 this system of drill-down reports allows all relevant
2 information to be compiled according to a constructed query
3 (date range, selected grammars, selected calls, etc.) for
4 purposes such as double-checking transcription accuracy,
5 application assessment, or insufficiently clear guides on
6 responses within a given grammar. Statistical and
7 heuristic analysis of the transcribed results compared to
8 the results of the speech recognizer in the context of the
9 grammar allow grammar authors and application programmers
10 to determine if the menu prompting options are sufficient
11 to guide a user through the menu as well as determining
12 whether the grammar and/or the pronunciation should be
13 tuned to be more consistent with typical menu use. For
14 example, if a certain pronunciation of a given word in a
15 grammar is consistently marked as mispronounced, the
16 grammar author might consider tuning the pronunciation
17 dictionary for the speech recognition software to include
18 that pronunciation of the word.

19 Figure 6 is a screen shot of a top-level drill-down
20 report according to one embodiment of the present
21 invention. Thus, web browser 200 navigates the browser
22 window to the address shown in address window 201 when
23 choosing the drill-down report feature of the present
24 transcription system. For example, a summary accuracy
25 report for a given date, shown in report screen 200(D), is
26 shown prior to each record to be transcribed. Data in
27 report screen 200(D) is organized into a table format,
28 wherein each column represents a type of top-level data
29 relevant for a top-level analysis of the accuracy of the
30 speech recognition system and each row represents a
31 different grammar. For example, columns 602-606 include
32 top-level data for the number of utterances 602,

1 classification of utterance 603, in-grammar performance
2 604, out-of-grammar performance 605, and overall
3 performance 606 data summaries for each corresponding
4 grammar in name of grammar column 601.

5 Specifically, in a telephone information service
6 having a menu which connects users to an airline of their
7 choice, the grammar for that menu includes the name of each
8 airline in the service. Thus, the grammar-in-use includes
9 airline names, such as "delta", "southwest", and "united".
10 The grammar-in-use additionally includes words in
11 applicable intrinsic grammars, such as "help" and "go
12 back". The Session.Airlines.Choice grammar, located in row
13 620 of accuracy report window 200(D), is the grammar for
14 such a telephone information service. As shown, 3000
15 utterances have been transcribed (row 620, column 602)
16 relating to the Session.Airlines.Choice grammar. These
17 utterances have been analyzed to provide the data present
18 in row 620, columns 603-606. Thus, of those 3000
19 utterances, 76.07% are in-grammar (column 603A) and 23.93%
20 are out-of-grammar (column 603B), where "out-of-grammar"
21 indicates that the utterance was not one of the valid words
22 within the Session.Airlines.Choice grammar used for the
23 telephone information service.

24 Of the 76.07% in-grammar utterances (column 603A), the
25 speech recognizer correctly interpreted 96.89% (column
26 604A), falsely accepted 0.66% (column 604B), and falsely
27 rejected 0.66% (column 604C). A false acceptance occurs
28 when the utterance is out-of-grammar, yet the speech
29 recognizer interprets the utterance as in-grammar. A false
30 rejection occurs when the utterance is in-grammar, yet the
31 speech recognizer interprets the utterance as out-of-
32 grammar. The comparison is made between the transcribed

1 utterance and the word recognized by the speech recognizer,
2 such that the percentage of correctly interpreted
3 utterances is equivalent to the number of in-grammar
4 utterances interpreted by the speech recognizer that match
5 the corresponding transcribed utterance divided by the in-
6 grammar number of utterances.

7 Of the 23.93% out-of-grammar utterances (column 603B),
8 the speech recognizer correctly rejected 26.46% (column
9 605A) and falsely accepted 73.54% (column 605B). The
10 overall performance of the speech recognizer for the
11 Session.Airlines.Choice grammar is described in column 604,
12 with the percentage of correct acceptances divided by all
13 utterances, and is equal to 73.70% (column 606A).

14 Each grammar in accuracy report screen 200(D) has
15 similar top-level information. Note that additional top-
16 level information may be added to accuracy report screen
17 200(D) by adding to the number of columns. Additional
18 information is available for this top-level by clicking on
19 the associated hyperlink. For example, the
20 Session.Airlines.Choice grammar is underlined. In a web-
21 based system, this underline (and typically an associated
22 color) indicates a hyperlink. In one embodiment, clicking
23 on the Session.Airlines.Choice grammar hyperlink navigates
24 web browser 200 to another web page displaying the valid
25 words in-grammar for the Session.Airlines.Choice grammar.
26 In another embodiment, clicking on the
27 Session.Airlines.Choice grammar hyperlink opens an
28 additional web browser in which the web browser screen
29 displays the valid words in-grammar for the
30 Session.Airlines.Choice grammar. Support data for the data
31 in columns 603-606 are similarly accessed.

Figure 7 is a screen shot of a first-level-down drill-down report according to one embodiment of the present invention. Clicking on the 2.45% false accepts for in-grammar performance (row 620, column 604B, of Figure 6) navigates web browser 200 to in-grammar false accepts screen 200(E). Thus, web browser 200 navigates the browser window to the address shown in address window 201 when choosing the 2.45% false accepts for in-grammar performance (row 620, column 604B, of Figure 6) in the present transcription system. Data in in-grammar false accepts screen 200(E) is organized into a file system format, wherein each row includes a folder icon (e.g. folder 701, which in one embodiment is itself a hyperlink), an number indicating the frequency of a particular type of false accept (for example number 702), the transcribed utterance (for example transcribed utterance 703), and the result of the speech recognizer (for example result 704). Key 720 describes the format for naming these in-grammar false accepts. Specifically, in row 710, the speech recognizer mistook the in-grammar utterance "help" (transcribed utterance 703) for the word "delta" (result 704) seven times (number 701). Similarly, in row 711, the speech recognizer mistook the in-grammar utterance "southwest" for the word "conquest" twice.

Note that the number of in-grammar utterances for the Session.Airlines.Choice grammar is the number of utterances (3000 in row 620, column 602, in Figure 6) multiplied by the percent of utterances in-grammar (76.07% in row 620, column 603A, in Figure 6), which is equivalent to 2286 in-grammar utterances. The number of false accepts of these in-grammar utterances is 2.45% (row 620, column 604B, Figure 6) multiplied by 2286 in-grammar utterances, which

1 is equivalent to 56 in-grammar false accepts. This number
2 of in-grammar false accepts is listed in line 721 of in-
3 grammar false accepts screen 200(E).

4 Additional information is available for this first-
5 level-down information by clicking on the associated folder
6 hyperlinks. Clicking on the folder 701 hyperlink (row 710)
7 opens a sub-list of the seven (number 702) in-grammar help-
8 delta false accepts. Specifically, in one embodiment,
9 clicking on the folder 701 hyperlink alters in-grammar
10 false accepts screen 200(E) to include hyperlinks to
11 ".wav", or "WAV format", files 801-807 as shown in Figure
12 8. Hyperlinks to .wav files 801-807 are indented under
13 folder 701 to show that they are the seven utterance
14 recordings of the in-grammar utterance "help" which were
15 recognized as "delta" by the speech recognizer. Support
16 data for the data in row 711 (and other rows) is similarly
17 accessed.

18 In one embodiment, clicking on a hyperlink to one of
19 .wav files 801-807 (e.g. .wav file 801) navigates web
20 browser 200 to another web page displaying the utterance,
21 result, and a sound tool for playing the utterance. In
22 another embodiment, clicking on a hyperlink to one of .wav
23 files 801-807 (e.g. .wav file 801) opens an additional web
24 browser in which the web browser screen displays the
25 utterance, result, and a sound tool for playing the
26 utterance.

27 Figure 9 is a screen shot of a third-level-down drill-
28 down report according to one embodiment of the present
29 invention. Clicking on the hyperlink for .wav file 801
30 navigates web browser 200 to wav file screen 200(F)
31 displaying transcribed utterance 703 (e.g. "help"), result
32 704 (e.g. "delta"), and a sound tool 310 for playing the

1 utterance audio. As described with respect to the
2 transcription process, audio control panel 310 allows the
3 utterance audio to be played, as well as other audio
4 operations to be performed, such as changing the volume and
5 pausing the replay of the recording.

6 In this way, both top-level data and low level data
7 can be easily displayed and quickly obtained. For example,
8 a specific sound file included in the performance analysis
9 of in-grammar false accepts can be accessed in three clicks
10 from the top-level description of performance.

11 12 Other Embodiments

13 In one embodiment, the transcription tools and
14 accuracy reports are made available as part of a zero-
15 footprint remotely hosted development environment. See, US
16 Patent Application Ser. No. 09/592,241, entitled "Method
17 and Apparatus for Zero-Footprint Application Development",
18 having inventors Jeff C. Kunins, et. al., filed 13 June
19 2000. In such configuration, the transcriptionist will
20 frequently be the application developer or her/his
21 authorized agent. Additionally, utterance access will be
22 limited to those utterances made within the developer's own
23 application(s). For example, if the application was
24 accessed by a user through "Shopping", "Bookstore", only
25 the utterances for grammars within the "Bookstore" menu
26 item would be available to the developer for transcription.

27 In one embodiment, the transcription and accuracy
28 tools are a separately paid for component of the zero-
29 footprint development environment. In another embodiment,
30 the developer can specifically request that the hosting
31 sites (e.g. the hosting site 101) record utterances for

1 her/his application(s). In some embodiments, there may be
2 a charge for this service.

3 In another embodiment, developers can request
4 transcription of a predetermined number of utterances,
5 e.g., 10,000, from the provider of the zero-footprint
6 development environment (or their affiliates, etc.) for a
7 cost. Then the developer can simply use the accuracy
8 reports without the need for her/him to perform the
9 transcriptions.

10 The embodiments described above are illustrative only
11 and not limiting. For example, in other embodiments of the
12 invention, additional steps such as secured login and data
13 encryption may be added to the transcription process.
14 Moreover, data may be displayed in any form that clearly
15 conveys meaningful information during report generation.
16 Other embodiments and modifications to the system and
17 method of the present invention will be apparent to those
18 skilled in the art. Therefore, the present invention is
19 limited only by the appended claims.